



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Single cell analyses and machine learning define hematopoietic progenitor and HSC-like cells derived from human PSCs

**Citation for published version:**

Fidanza, A, Stumpf, PS, Ramachandran, P, Tamagno, S, Babbie, A, Lopez Yrigoyen, M, Taylor, H, Easterbrook, J, Henderson, B, Axton, R, Henderson, NC, Medvinsky, A, Ottersbach, K, Romanò, N & Forrester, L 2020, 'Single cell analyses and machine learning define hematopoietic progenitor and HSC-like cells derived from human PSCs', *Blood*, vol. 136, no. 25, pp. 2893–2904.  
<https://doi.org/10.1182/blood.2020006229>

**Digital Object Identifier (DOI):**

[10.1182/blood.2020006229](https://doi.org/10.1182/blood.2020006229)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

Blood

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



## Single cell analyses and machine learning define hematopoietic progenitor and HSC-like cells derived from human PSCs.

Tracking no: BLD-2020-006229R1

Antonella Fidanza (University of Edinburgh, United Kingdom) Patrick Stumpf (University of Southampton, United Kingdom) Prakash Ramachandran (University of Edinburgh, United Kingdom) Sara Tamagno (University of Edinburgh, United Kingdom) Ann Babbie (Centre for Integrative Systems Biology and Bioinformatics, Department of Life Sciences, Imperial College London, United Kingdom) Martha Lopez-Yrigoyen (University of Edinburgh, United Kingdom) Alice Taylor (University of Edinburgh, ) Jennifer Easterbrook (University of Edinburgh, United Kingdom) Beth Henderson (University of Edinburgh, United Kingdom) Richard Axton (University of Edinburgh, United Kingdom) Neil Henderson (University of Edinburgh, United Kingdom) Alexander Medvinsky (The University of Edinburgh, United Kingdom) Katrin Ottersbach (University of Edinburgh, United Kingdom) Nicola Romano (University of Edinburgh, United Kingdom) Lesley Forrester (University of Edinburgh, United Kingdom)

### Abstract:

Haematopoietic stem and progenitor cells (HSPCs) develop through distinct waves at various anatomical sites during embryonic development. The *in vitro* differentiation of human pluripotent stem cells (hPSCs) is able to recapitulate some of these processes but it has proven difficult to generate functional haematopoietic stem cells (HSCs). To define the dynamics and heterogeneity of HSPCs that can be generated *in vitro* from hPSCs, we exploited single cell RNA sequencing (scRNAseq) in combination with single cell protein expression analysis. Bioinformatics analyses and functional validation defined the transcriptomes of naïve progenitors as well as erythroid, megakaryocyte and leukocyte-committed progenitors and we identified CD44, CD326, ICAM2/CD9 and CD18 as novel markers of these progenitors, respectively. Using an artificial neural network (ANN), that we trained on a scRNAseq derived from human fetal liver, we were able to identify a wide range of hPSCs-derived HPSC phenotypes, including a small group classified as HSCs. This transient HSC-like population reduced as differentiation proceeded and was completely missing in the dataset that had been generated using cells selected on the basis of CD43 expression. By comparing the single cell transcriptome of *in vitro*-generated HSC-like cells with those generated within the fetal liver we identified transcription factors and molecular pathways that can be targeted with the aim of improving HSC differentiation *in vitro*.

**Conflict of interest:** No COI declared

**COI notes:**

**Preprint server:** Yes; BioRxiv <https://doi.org/10.1101/602565>

**Author contributions and disclosures:** AF, designed and performed research, analyzed the data and wrote the manuscript. AF, PS, AB and NR performed bioinformatics analysis. PR, ST, MLY, AHT, JE, BH, RA performed research. LMF designed the experiment, analyzed data and wrote the manuscript. NH, AM, KO, and NR provided intellectual input and final approval of the manuscript.

**Non-author contributions and disclosures:** No;

**Agreement to Share Publication-Related Data and Data Sharing Statement:** We have created a webpage where the data can be freely browsed, plots can be generated and exported, and full datasets can be downloaded. The link is provided in the manuscript.

**Clinical trial registration information (if any):**

# Single cell multimodal analyses and machine learning define haematopoietic progenitor and HSC-like cells derived in vitro from human pluripotent stem cells.

Antonella Fidanza<sup>1\*</sup>, Patrick S Stumpf<sup>2</sup>, Prakash Ramachandran<sup>3</sup>, Sara Tamagno<sup>1</sup>, Ann Babbie<sup>4</sup>, Martha Lopez-Yrigoyen<sup>1</sup>, A. Helen Taylor<sup>1</sup>, Jennifer Easterbrook<sup>1</sup>, Beth Henderson<sup>3</sup>, Richard Axton<sup>1</sup>, Neil C. Henderson<sup>3</sup>, Alexander Medvinsky<sup>1</sup>, Katrin Ottersbach<sup>1</sup>, Nicola Romanò<sup>5</sup>, Lesley M. Forrester<sup>1\*</sup>.

1 - Centre for Regenerative Medicine, University of Edinburgh, Edinburgh, UK

2 - Joint Research Center for Computational Biomedicine, Uniklinik RWTH Aachen, Aachen, Germany

3 - Centre for Inflammation Research, University of Edinburgh, Edinburgh, UK

4 - Centre for Integrative Systems Biology and Bioinformatics, Department of Life Sciences, Imperial College London, London, UK

5 - Centre for Discovery Brain Sciences, University of Edinburgh, Edinburgh, UK

\*correspondence to [afidanza@ed.ac.uk](mailto:afidanza@ed.ac.uk) and [L.Forrester@ed.ac.uk](mailto:L.Forrester@ed.ac.uk)

**Key point 1** - Single-cell and CITE-seq profiling of human HSPCs derived in vitro from pluripotent stem cells browsable at <http://188.166.158.65/scRNAseqHPC/>

**Key point 2** – Artificial Neural Network identifies HSC-like cells derived in vitro from hPSCs.

## Abstract

Haematopoietic stem and progenitor cells (HSPCs) develop through distinct waves at various anatomical sites during embryonic development. The in vitro differentiation of human pluripotent stem cells (hPSCs) is able to recapitulate some of these processes, however, it has proven difficult to generate functional haematopoietic stem cells (HSCs). To define the dynamics and heterogeneity of HSPCs that can be generated in vitro from hPSCs, we exploited single cell RNA sequencing (scRNAseq) in combination with single cell protein expression analysis. Bioinformatics analyses and functional validation defined the transcriptomes of naïve progenitors as well as erythroid, megakaryocyte and leukocyte-committed progenitors and we identified CD44, CD326, ICAM2/CD9 and CD18

as markers of these progenitors, respectively. Using an artificial neural network (ANN), that we trained on a scRNAseq derived from human fetal liver, we were able to identify a wide range of hPSCs-derived HPSC phenotypes, including a small group classified as HSCs. This transient HSC-like population decreased as differentiation proceeded and was completely missing in the dataset that had been generated using cells selected on the basis of CD43 expression. By comparing the single cell transcriptome of in vitro-generated HSC-like cells with those generated within the fetal liver we identified transcription factors and molecular pathways that can be exploited in the future to improve the in vitro production of HSCs.

## Introduction

Human pluripotent stem cells (hPSCs) can be differentiated in vitro into various haematopoietic cell types, providing both a model for basic research studies and a source of clinically relevant cells<sup>1</sup>. During embryonic development, two waves of restricted haematopoietic progenitors arise in the extraembryonic tissues of the yolk sac, before emergence of haematopoietic stem cells (HSCs) in the embryo proper<sup>2</sup>. In the mouse embryo, at E7.25 the first “primitive” wave gives rise to erythrocytes, megakaryocytes and macrophages<sup>3,4</sup>, after at E8.25, the second wave, also known as the first “definitive” progenitors, consists of erythro-myeloid progenitors (EMPs) that can be distinguished from the primitive progenitors by their potential to generate granulocytes<sup>5</sup>. Furthermore, the monocytes that emerge from EMPs provide the embryo with tissue resident macrophage; the first life-long lasting population of immune cells<sup>6–8</sup>. Intraembryonic hematopoiesis is established during E10.5-E11.5 in the aorta-gonad-mesonephros (AGM) region with the emergence of HSCs, that will sustain the lifespan production of all blood lineages, also upon transplantation<sup>9</sup>. A number of studies have indicated that human haematopoietic development follows a comparable process<sup>10–13</sup> but for ethical reasons it has proven difficult to gain a clear insight into the lineage potential and hierarchical relationships between early human haematopoietic progenitors. The dynamic nature and the heterogeneity of haematopoietic progenitor populations that arise during development poses additional confounders to the identification of molecular mechanisms associated with their formation and function.

To gain insight into the transcriptome of developing human haematopoietic progenitors, we performed in-depth characterization of haematopoietic progenitors derived from

hPSCs. Single-cell expression profiles of hPSCs-derived haematopoietic cells have been reported but these previous studies either used a limited number of cells or used biased approaches in their isolation and so failed to depict their trajectory of differentiation<sup>14–16</sup>. This significantly impacted on the ability to resolve the complex heterogeneity of the progenitor pool, to identify the hierarchical relationship between subpopulations and to compare the transcriptome of hPSCs-derived progenitors to their in vivo counterparts. To address these issues, we generated scRNAseq data sets of human hPSCs-derived haematopoietic progenitors. Lineage trajectories predicted in silico were validated by functional assays of sorted cells and the specificity of our marker repertoire was confirmed using a CITE-seq<sup>17</sup> strategy. Furthermore, to annotate the hPSCs-derived progeny in an unbiased manner, we employed machine learning and trained an artificial neural network (ANN) to recognize the single-cell gene expression profiles of human fetal liver cells. This trained ANN was subsequently used to predict the identities of hPSCs-derived cells. The ANN thereby provides a mapping between in vivo and in vitro hematopoiesis and revealed a subset of hPSCs-derived cells that closely resembles HSCs in the foetal liver. Finally, by comparing that transcriptome of in vitro and in vivo-generated HSCs we identified genetic pathways that can be exploited to improve HSCs production in vitro from hPSCs.

## Methods

Methods are available as supplementary methods.

## Results

### Single cell RNA sequencing of iPSCs-derived haematopoietic progenitor cells reveals the transcriptome of naïve and lineage committed progenitors.

To resolve the heterogeneity of in vitro generated hPSCs-derived haematopoietic progenitors we designed a minimal membrane marker strategy that allows to broadly isolate hPSCs-derived haematopoietic progenitors. This marker strategy was validated using two hPSC reporter lines, RUNX1C-GFP and KLF1-mCherry, CFU-C assays of sorted cell populations and gene expression profiling (Supplementary Figure S1). These data, together with previous reports<sup>18–22</sup>, supported our rationale that the isolation of CD235a<sup>+</sup>CD43<sup>+</sup> cells would enrich for HSPCs and exclude cells derived from the primitive wave (Supplementary Figure S1). We anticipated that the CD235a<sup>+</sup>CD43<sup>+</sup> compartment

would also comprise the early stages of lineage commitment, capturing the downstream hierarchy of early human progenitors.

CD235a<sup>+</sup>CD43<sup>+</sup> suspension cells from two independent replicate cultures at day 13 of differentiation were sorted by FACS and subjected to scRNAseq and data analyses (Figure 1A). After quality control and clustering<sup>23</sup> we obtained the transcriptome of 11420 cells (Supplementary Figure 2A-C) belonging to eight clusters of cells (Figure 1B). Although the two replicates did not show obvious differences (Supplementary Figure 2C), any potential batch effect was regressed out prior to integration. We assigned cell identities based on the expression of known markers and identified additional markers from the dataset that were cluster specific (Figure 1C-D). Clusters containing more immature, uncommitted progenitors were identified by their expression of progenitor-associated genes such as *KIT* and *GATA2* and their lack of expression of genes associated with specific cell lineages and were thus annotated as naïve populations (Figure 1D, Supplementary Figure 2D). Clusters that displayed expression of lineage markers were annotated as lineage committed progenitors (Figure 1B-D, Supplementary Figure 2D), including clusters of cells committed towards the megakaryocyte (*GP9* and *PF4*), erythroid (*GYP*A and *KLF1*) and granulocyte (*AZU1* and *PRNT3*) lineages (Figure 1D). Markers for each of the cell clusters were identified by differential gene expression analysis, further supporting the identities assigned to each cluster (Figure 1C, Supplementary Table 1).

## **Trajectory analyses reveal the hierarchy of in vitro derived haematopoietic progenitors.**

To study the hierarchical relationship between cell populations, we performed trajectory analysis using different methods including diffusion analysis<sup>24</sup> using Seurat R package<sup>23</sup> and pseudotemporal ordering, using Monocle R package<sup>25</sup> and Partition-based graph abstraction (PAGA)<sup>26</sup>. Diffusion analysis identified a central core from which three distinct trajectories emerged (Figure 2A). The central core corresponded to cells that we had annotated as naïve progenitors (Figure 2A-B). Branches comprised cells expressing genes associated with specific lineages, annotated as Erythroid (Ery)-, Megakaryocyte (Mega)- and Granulocyte (Granulo)-committed lineages. Comparable trajectories were observed using pseudotemporal ordering with PAGA and Monocle (Figure 2C-D). The PAGA analysis showed that the naïve cells were highly connected to erythroid, megakaryocyte

and granulocyte committed cells (Figure 2C). Pseudotime reconstruction of the hierarchy showed that cells annotated as naïve 1 were located at the top of the hierarchy and appeared to progress to naïve 2 cells before entering branches containing lineage committed cells (Figure 2D-E). Lineage commitment was also inferred from the expression of lineage-associated transcription factors that were filtered from the cluster specific marker genes according to their GO annotation (Figure 2F, **Supplementary Figure 2E**). For example, erythroid committed clusters demonstrated expression of both *KLF1* and *MYC*, with the latter decreasing in Ery 2 compared to Ery1, in keeping with their position within the hierarchy (Figure 2E, F). Within the megakaryocyte-committed clusters 1 and 2 we observed the expression of *GATA1*, *TAL1* and *FLI1* a cocktail of genes recently used for hPSCs forward programming to megakaryocytes (Figure 2F)<sup>27</sup>. Granulocyte-committed cells were represented by a separate branch and demonstrated the expression of *CEBP-D*, *CEBP-B*, *CEBP-A* and *CEBP-E* (Figure 2F). We then focused our attention on the transcription factors expressed by the naïve progenitor clusters and noted a high level of expression of *LMO4* and *ID2*, as well as *GATA2* which is known to be expressed in HSPCs (Figure 2F). We then inferred their role in the gene network using a network inference algorithm (Partial Information Decomposition and Context, PIDC)<sup>28,29</sup>. Single cell transcriptomic data offers the potential to observe dependencies between the expression profiles of pairs of genes, that if co-regulated, are expected to change in a coordinated fashion. Genes with highest statistical dependencies are connected by edges that altogether define the network<sup>28,29</sup>. Many of the transcription factors previously detected to be highly expressed within the naïve cell populations such as *ID2*, *ID4* and *LMO4*, occupy nodes within this large network (Supplementary Figure 3A-B). This strategy corroborates the importance of the identified transcription factors as functional elements within the single cell gene network.

### **CD44 membrane expression marks human clonogenic haematopoietic progenitors.**

To experimentally validate the results of our trajectory analyses experimentally, we set out to assess the haematopoietic potential of the naïve progenitor populations. We defined a prospective sorting strategy using genes encoding the cell surface markers *CD33*, *CD44*, and *ITGB2* (also known as CD18) that were enriched within the naïve progenitors' clusters (Figure 3A). *CD33* was expressed by both naïve 1 and naïve 2 progenitors whereas *CD44* and *CD18* expression appeared higher in the naïve 1 population (Figure 3A). We



166 fractionated CD235a<sup>-</sup>CD43<sup>+</sup>CD33<sup>+</sup> cells and identified subpopulations as naïve 1A  
 167 (CD44<sup>+</sup>CD18<sup>-</sup>), naïve 1B (CD44<sup>+</sup>CD18<sup>+</sup>) and naïve 2 (CD44<sup>-</sup>CD18<sup>-</sup>) (Figure 3B). Trajectory  
 168 analysis predicted that naïve 1 cells were at the top of the hierarchy and gave rise to the  
 169 naïve 2 cells prior to lineage commitment (Figure 2D-E). To test this in silico prediction, we  
 170 used a chimeric co-culture system using the Zeiss Green (ZsG) reporter (Figure 3C). This  
 171 approach allowed us to sort, for example, ZsG-labelled naïve 1 cells, then track their ZsG  
 172 progeny after being placed back in the complex differentiation environment. We  
 173 synchronously differentiated the ZsG-iPSC line, constitutively expressing the fluorescent  
 174 reporter<sup>30</sup>, and the parental iPSC line. To verify the progressions of naïve 1 to naïve 2 and,  
 175 naïve 2 to lineage committed cells, we sorted naïve 1 (CD33<sup>+</sup>CD44<sup>-</sup>CD18<sup>-</sup>) or naïve 2  
 176 (CD33<sup>+</sup>CD44<sup>+</sup>CD18<sup>-/+</sup>) cells from ZsG-iPSCs at day 10 and co-cultured these with the  
 177 synchronized differentiating parental cells for a further 3 days. As predicted from the  
 178 trajectory analysis, the naïve 1 cell population was able to generate ZsG-expressing naïve  
 179 2 cells. We also noted that the naïve 1 cells retained their immunophenotype, indicating  
 180 some self-renewal capacity (Figure 3D). Interestingly, naïve 2 cells demonstrated some  
 181 potential to acquire CD44 and CD18, markers of naïve 1 cells (Figure 3D), suggesting  
 182 fluidity between these states. As predicted by our trajectory analyses (Figure 2D-E), naïve  
 183 2 cells acquired the ability to generate more mature cells including erythroid cells  
 184 (CD235a<sup>+</sup>), megakaryocytes (CD41<sup>+</sup>) and macrophages (25F9<sup>+</sup>) (Supplementary Figure  
 185 3C). We compared the colony forming capacity of naïve 1 and 2 progenitors present at day  
 186 10 and day 13. When plated in clonogenic CFU-C assays, CD44<sup>+</sup> naïve 1 cells formed  
 187 CFU-C colonies but virtually no colonies were generated by naïve 2 cells at either time  
 188 point (Figure 3E-F). These data support the proposed hierarchy and indicate that CD44  
 189 expression alone resolves colony forming cells. Our chimeric co-culture system was able  
 190 therefore to assess the lineage output that could not be assessed solely by CFU-C assays.  
 191 We observed that naïve progenitors expressed high levels of ID genes (Figure 2F), and  
 192 that they were identified as nodes within the gene network (Supplementary Figure 3A). As  
 193 ID genes are targets of BMP signaling, we predicted that these naïve cells would be  
 194 responsive to BMP stimulation. We added BMP4 to differentiation culture from day 10,  
 195 when both naïve 1 and 2 were present and then assessed the proportion of these cells 3  
 196 days later. In presence of BMP4, we observed a 25% and 59% expansion of naïve 1 and 2  
 197 cells respectively (Supplementary Figure 3E). In this experiment we used both hESCs and  
 198 hiPSCs and showed that naïve progenitors are present at a comparable frequency in both



hESCs and hiPSCs (Supplementary Figure 3D), and that naïve progenitors derived from both lines responded to BMP stimulation in a comparable manner (Supplementary Figure 3D-E). Thus this experiment not only identified an important functional signaling pathway but also confirmed that the markers we used to define naïve progenitors, and their response to BMP signaling, are not PSC line specific.

To assess whether the naïve cell populations identified using our unique sorting strategy showed features of definitive haematopoietic progenitors, we assessed the expression of the RUNX1C-GFP reporter. We observed RUNX1C-GFP expression in both cell types, with a higher proportion of RUNX1C<sup>+</sup> cells in the naïve 1 compared to naïve 2 population (Figure 3G). Definitive HSPCs are generated via endothelial to hematopoietic transition (EHT) during embryonic development<sup>31,32</sup> so, we would expect comparable hPSCs-derived progenitors to have hallmarks of their endothelial origin. Here we demonstrate that naïve CD44<sup>+</sup> cells generated in vitro from hPSCs co-expressed CD34 and the endothelial marker CD144 (also known as VeCad) with expression being higher at day 10, when the majority of naïve progenitors were present (Supplementary Figure 3L). This endothelial signature of naïve progenitors, together with their lineage potential reflects their definitive features. To confirm that CD44 expression was associated with HSPCs in vivo we demonstrated its co-localization with CD45 and CD144 in the mouse yolk sac and AGM region (Supplementary figure 3F-J). At E10.5 in the yolk sac, CD44 was expressed on endothelial cells in a bimodal fashion, with vessels expressing low and high levels, the latter being associated with bright clusters of haematopoietic cells (Supplementary figure 3G). By flow cytometry, we observed that by E11, all CD45<sup>+</sup> cells and a proportion of CD144<sup>+</sup> cells were within the CD44<sup>high</sup> population (Supplementary figure 3F). Within the embryo proper, CD44 was expressed on the membrane of endothelial cells within the dorsal aorta, whereas venous endothelial layers were CD44 negative (Supplementary figure 3H-I). CD44 was also co-expressed with CD45<sup>+</sup> in the AGM region (Supplementary figure 3H-J). Furthermore, expression of LMO4 in CD44<sup>+</sup> cells within the AGM region is in keeping with its high level of expression in naïve hPSCs-derived HSPCs (Figure 2F) and supports its identification as a novel haematopoietic transcription factor. These data suggest that CD44 is expressed on haemogenic endothelial cells and it is retained on emerging haematopoietic cells in vivo, similar to what we have observed during the in vitro differentiation of human progenitors (Supplementary Figure 3I).

## Identification of membrane markers of lineage committed progenitors

We next selected membrane markers that we predicted could be used for the isolation of lineage committed progenitors. Erythroid-primed clusters 1 and 2 both showed expression of *MYC* (Figure 2F) and *EPCAM* (also known as CD326) (Supplementary Figure 4A), indicative of early committed erythroid cells<sup>33,34</sup>. We confirmed that CD326 was expressed in the majority of CD235a<sup>+</sup> cells at day 13 of iPSC differentiation but interestingly, we noted a small number of CD326<sup>+</sup>CD235a<sup>-</sup> (Supplementary Figure 4B), suggesting that CD326 might be marking commitment to the erythroid lineage prior to CD235a acquisition. To test this, we assessed the expression dynamics of these markers during the erythroid differentiation of umbilical cord blood CD34<sup>+</sup> (UCB34<sup>+</sup>) cells. At day 10 of differentiation, CD326 was expressed in CD235a<sup>-/low</sup> cells but not in CD235a<sup>high</sup> cells, the latter corresponding to more mature erythroid cells (Supplementary Figure 4B). CD326 was not expressed in cells at day 18 of the differentiation protocol (when the majority of cells are mature CD235a<sup>+</sup> cells) nor in the mature erythrocytes found in adult peripheral blood (Supplementary Figure 4B). Taken together these data suggest that CD326 marks early erythroid progenitors in both hiPSC, fetal and adult derived cells. We also noted the expression of *HBG1*, *HBG2*, *HBA1*, and *HBA2*, subunits of fetal hemoglobin, indicative of erythroid cells derived from definitive hematopoiesis (Supplementary Figure 4C).

Three clusters with megakaryocyte and platelet signatures (Mega-primed 1, 2 and 3) were predicted by expression of *ITGA2B* (CD41), *GP9*, *PF4* (Figure 1C-D and Supplementary Table 1). *ICAM2* was expressed at higher level in cluster Mega-primed 3 (Supplementary Figure 4D), as for CD9, known to increase along megakaryocytes differentiation<sup>35</sup>. *ICAM2* and CD9 co-expression was confirmed by flow cytometry (Supplementary Figure 4D). We observed a population of CD41<sup>+</sup>CD9<sup>+</sup>ICAM2<sup>+</sup> cells, with around 85% of the CD41<sup>+</sup>CD42a<sup>+</sup> (Supplementary Figure 4E), that did not detect polyploidy, supporting their immature status (Supplementary Figure 4F-G).

Granulocyte-committed clusters were predicted by expression of markers such as *MPO*, *AZU1*, *RNASE2* and *ITGB2* which encodes the membrane marker CD18, subunit of different leukocytes marker such as CD11a-d, Mac-1 and LFA-1 (Figure 1C, Supplementary Table 1). We sorted CD235a<sup>-</sup>CD43<sup>+</sup>CD33<sup>+</sup>CD44<sup>-</sup>CD18<sup>+</sup> cells and confirmed the phenotype of granulocytes and monocytes based on their nuclear morphology (Supplementary Figure 4H). Further clustering revealed three sub-clusters corresponding to eosinophil, neutrophils and monocytes lineages (Supplementary Figure

4I-L). Noteworthy, *RUNX3* expression was specifically associated with the monocyte subcluster (Supplementary Figure 4J) previously reported to be important for zebrafish stem cells and macrophages<sup>36</sup>, and to be expressed by mouse tissue resident macrophages of the skin<sup>37</sup>.

In summary, we showed that naïve progenitors give rise also to committed progenitors characterized by features of granulocytes and monocyte, cell types that emerge exclusively in the definitive waves<sup>5</sup>.

### **CITE-seq analysis of human iPSC-derived haematopoietic progenitors.**

To further study the temporal emergence of the progenitor populations and their associated markers, we carried out CITE-seq analysis whereby single cell membrane marker expression can be directly correlated with the relative transcriptome<sup>17</sup>. To ensure that we sampled even the rarest progenitor cell population we extended the CITE-seq analysis to adherent cells and included an earlier time point (day 10) in addition to day 13. Again, to exclude primitive erythroid cells, we selected CD235a-negative suspension cells but, in this experiment, we included and enriched for CD43<sup>+</sup> cells that had been excluded in our initial study. (Supplementary figure 5B). We expected early progenitors to express CD31 and to potentially remain part of hematopoietic clusters within the adherent fraction of the culture and so we FAC-sorted the adherent cells into CD31<sup>+</sup> and CD31<sup>+</sup> fractions.

Cells were labeled with oligonucleotide tagged antibody specific for the membrane markers that we identified in our initial experiment (ADT\_CD18, ADT\_CD33, ADT\_CD41, ADT\_CD44, ADT\_CD102, ADT\_CD326; ADT: Antibody-Derived Tag) as well as other markers of endothelial and early committed hematopoietic cells (ADT\_CD144) and of macrophages (ADT\_CD163). To test the specificity of the membrane marker repertoire previously identified on the suspension cells, we subset the two libraries corresponding to suspension cells collected at day 10 and 13 (Figure 4, Supplementary Figure 5B-C). After multidimension reduction and clustering analysis we identified a naïve progenitor population (Figure 4A), comparable to our first sequencing experiment (Figure 2A). These naïve progenitors exhibited erythroid (Ery), megakaryocyte (Mega), and granulocyte and monocytes (Gra-Mo) lineage potential, with increased lineage commitment at day 13 compared to day 10 (Figure 4B); in line with the expression pattern of genes associated with naïve and committed stages in these days (Supplementary Figure 5D). Analysis of the ADTs showed that each marker was expressed in the expected cluster (Figure 4C) thus

supporting them as markers for defined progenitors. To further explore the power of the ADT approach, we performed multidimension reduction using ADTs as the only input dimensions and proved that ADT data alone identified remarkably similar clusters (Figure 4D-E), strongly correlated with the clusters derived from the entire transcriptome (Figure 4F). Taken together, the CITE-seq approach confirms that the markers identified from our first scRNAseq analysis define the hierarchy of human developmental hematopoiesis in vitro with high specificity.

### **Comparison of in vitro generated progenitors with in vivo produced cells.**

The use of human PSCs as a renewable source of hematopoietic cell types faces major challenges relating to, for example, the inefficient repopulation capacity of progenitor cells and the incomplete maturation of differentiated cell types. To identify the underlying molecular basis associated with these deficiencies in hPSC-derived cells, we compared our dataset to a human fetal liver dataset which contains the complete hematopoietic hierarchy from long-term reconstituting HSCs to differentiated cell types.

To assess how hPSCs-derived naïve and lineage-committed progenitors compared to their equivalent counterpart generated in vivo, we assessed the expression of selected genes identified to distinguish the various cell types detected in the human fetal liver<sup>38</sup> (Figure 4G). An initial analysis of marker genes of lineage commitment in the developing embryo revealed that these markers are remarkably powerful for discriminating the equivalent in vitro cell types identified in our in vitro study (Figure 4G, Supplementary Table 1).

Interestingly, *SPINK2*, a newly reported marker of fetal HSC/MPP<sup>38</sup>, was also expressed specifically by our naïve progenitor cells (Figure 4G), together with CD34 (Supplementary figure 3L). These specific similarities observed between in vitro and in vivo developing hematopoietic progenitor cells led us to investigate in a more comprehensive manner the phenotype of cell types that are produced in vitro and how well these in vitro derived cells reflect the corresponding cell types during in vivo development. Therefore, we used the same published human fetal liver scRNAseq data as a reference, firstly, to identify in vitro derived cells with gene expression signatures of human fetal liver hematopoietic cells and, secondly, to quantify the similarity to their corresponding transcriptomes. To address the

first question, we employed machine learning to transfer labels from the fetal liver reference data to our in vitro-derived blood cells (Figure 5A). This approach enabled a much broader and unbiased identification of cell types compared to inference based purely

on marker genes. We followed our recently developed strategy<sup>39</sup> and trained an artificial neural network (ANN)<sup>39</sup> to recognize single-cell gene expression profiles of human foetal liver cells that were sampled at a time in development at which the liver is the main site of blood cell formation<sup>38</sup>. Briefly, this ANN is trained using the expression data of 3,479 genes and 145,725 cells from fetal liver as an input<sup>38</sup>. From these labelled data, the ANN learns to predict, from which of the 28 different fetal liver cell types a particular gene expression pattern originates. Once trained, the ANN is given previously unseen test data from in vitro derived cells as an input in order to annotate these data with human fetal liver cell labels. Since this approach considers 3,479 genes, it enabled a more comprehensive identification of cell types based on similarities in global gene expression patterns rather than specific marker genes.

The ANN was able to identify cell types within the source domain (the fetal liver data) with high accuracy as shown by the performance metrics obtained from 5-fold cross-validation (Supplementary figure 6A-B). The trained ANN was subsequently applied to the target domain (in vitro) to test if the hPSCs-derived cells were similar to those present in the foetal liver, in which case the label of that specific in vivo cell would be transferred. The ANN was able to assign labels to 92% of in vitro produced cells into various cell types present in vivo (Supplementary figure 6 C-D), most notably, a small population was labeled as HSC/MPP. This indicates that the global gene expression pattern of a subset of the in vitro derived cells is very similar to HSC/MPPs from the in vivo reference data in fetal liver. To quantify precisely how similar these in vitro derived HSC/MPPs are to their in vivo counterparts, we calculated the average pairwise Euclidean distance between HSC/MPPs, using the human fetal liver as a reference. This analysis indicates that fetal liver HSC/MPPs are, on average, only marginally more similar to one another as they are to iPSC derived HSC/MPPs (Supplementary Figure 8A). In summary, this analysis indicates that the in vitro derived HSC/MPPs closely, yet not perfectly, reflect the gene expression patterns of their in vivo counterparts. Using the ANN we also observed that the relative abundance of the predicted HSC/MPP population decreased with time by day 13 (Figure 5B), whereas, the relative abundance of committed cells increased over this time as expected (Supplementary Figure 6E). When we applied the same ANN strategy to our first data set, that was generated from day 13 progenitors that were selected on the basis of CD43 expression, no HSC/MPP were detected (Figure 5C). This is consistent with our observation that this transient HSC/MPP population is present in higher numbers earlier at

day 10, when they are almost equally distributed in the adherent CD31+ and suspension  
 CD235a- compartment (Supplementary Figure 6F). We looked for marker genes that  
 defines this predicted HSC/MPP cell population in vitro and looked specifically for  
 membrane markers according to their GO annotation (Supplementary Table 1). Together  
 with expected markers such as *CD34*, *CD44* and *CD33*, we also detected *CD132*, *CD52*,  
*CD180* and *IL3RA* and many others that will allow to design a prospective sorting strategy  
 to isolate this specific population. We then subset the in vivo and in vitro HSC/MPP and  
 integrated the two datasets (Figure 5D). The integrated data allowed for direct comparison  
 of their transcriptome and identified 54 differentially expressed genes (Supplementary  
 Table 1), all of which were lower in HSC/MPP produced in vitro compared to those  
 generated in vivo. GO analysis of these genes identified enrichment for KEGG signaling  
 pathways such as NOD-like receptor, IL-17, NF-Kappa B and HIF-1 (Supplementary Table  
 1). We also identified 6 genes encoding transcription factors: *EGR1*, *ZFP36L1*, *NR4A1*,  
*FOS*, *JUN* and *JUNB* (Figure 5E). Interestingly, the EGR1 binding site was enriched,  
 amongst others, in the upstream region of the differentially expressed genes (Figure 5F),  
 suggesting an important regulatory role of EGR1.

We also compared the predicted HSC/MPP derived from hPSCs to hematopoietic  
 progenitors isolated from different sites of hematopoiesis in the developing embryo  
 including to fetal liver HSC/MPP<sup>38</sup>, yolk sac MPP<sup>38</sup> that were collected at Carnegie stages  
 5 to 14, and AGM<sup>40</sup> sorted progenitors (CD34+CD45+CD235a-) collected at Carnegie  
 stage 15, around the time of early HSC emergence (Supplementary Figure 7A-B). Whole  
 transcriptome comparison, followed by KEGG pathway analysis, showed that in vitro  
 HSC/MPP cells are marked by genes associated with oxidative phosphorylation  
 (Supplementary Table 1), indicating metabolic differences between in vitro and in vivo  
 produced progenitors. Hypoxic conditions characterize mammalian embryo  
 development<sup>41</sup>, and more specifically the development of the hematopoietic system,  
 where hypoxia has been detected in hematopoietic clusters in the AGM region, and in the  
 fetal liver<sup>42</sup>. The hematopoietic progenitors derived from hPSCs were instead  
 differentiated in normoxic conditions which could explain their different metabolic profile.  
 The fetal liver cells were marked by HLA genes and consequently KEGG pathways  
 associated with antigen presentation and T-cell development (Supplementary Table 1).  
 The AGM dataset displayed high expression levels of genes associated with Notch  
 pathway, such as *HES1*, *NOTCH1*, *NOTCH2*, *JAG1* and *JAG2*. This is in line with the



developmental stage at which they were collected when the Notch pathway is orchestrating the HSC emerge<sup>43</sup>. Within the markers of yolk sac progenitors, we detected genes related to early hematopoietic development. *FRZB*, mesodermal cell marker, and *HBE1*, marker of primitive hematopoiesis, were listed in the top 10 differentially expressed genes: this underlines the early developmental features of yolk sac progenitors. Finally, we noted also that *SPINK1* was identified as marker for YS progenitors. While *SPINK2*, identified here and by others as marker of progenitor cells<sup>38,40</sup> was expressed by progenitors from all the tissues, *SPINK1* was detected exclusively in the YS (Supplementary Figure 7X), suggesting that this gene could discriminate extraembryonic from intraembryonic hematopoiesis.

Finally, we compared lineage committed cells identified by the ANN, in our in vitro dataset, to their in vivo counterpart from fetal liver, to identify genes that can be used as targets to improve the production in vitro of differentiated blood cell types. We listed the differentially expressed genes between in vitro and fetal liver cells and identified the transcription factors within the list (Supplementary Figure 8, Supplementary Table 1). Particularly interesting, late erythroid cells in vitro show high level of *PCLG2*, phospholipase C gamma 2, able to control intracellular calcium via production of IP3, inositol triphosphate. Intracellular calcium peaks just before enucleation, prior nuclei extrusion in the orthochromatic erythroblasts<sup>44</sup>. Erythroid cells derived from hPSCs are characterized by a general inefficient enucleation<sup>45,46</sup>, independently of their primitive or definitive origin (Supplementary Figure 1J) and this could be related to their intracellular calcium control.

In summary, we have identified a rare population of HSC/MPP-like cells in vitro that emerge early during differentiation of hPSCs and that display broadly similar gene expression patterns when compared to HSCs in development. However subtle differences are also apparent and a more detailed study of these differences could explain the known deficiencies of PSC-derived cells and ultimately be exploited to improve their therapeutic use. Our novel approach combines scRNAseq and machine learning to help identify candidate genes that may improve the production of HSCs and mature lineage cells from pluripotent stem cells in vitro, by closely recapitulating in vivo hematopoiesis.

## Discussion



430 We described the single cell transcriptome and membrane markers of naïve hematopoietic  
 431 progenitors and their lineage committed descendants derived in vitro from human  
 432 pluripotent stem cells. The repertoire of membrane markers proved to be remarkably  
 433 accurate in capturing the different states prior to and after lineage commitment.  
 434 We identified a population of naïve progenitors situated at the top of the differentiation  
 435 hierarchy, marked by CD44, a protein involved in the hematopoietic transition of the  
 436 hemogenic endothelium in the mouse AGM region<sup>47</sup>. We validated their lineage potential  
 437 employing a chimeric culture system, where isolated naïve progenitors, marked by Zeiss-  
 438 Green expression, demonstrated overlapping lineage output to that predicted in silico.  
 439 We also observed that progenitors are capable of moving between the naïve states, as  
 440 well as progressing into committed states. This is in keeping with many other scRNAseq  
 441 and proteomic studies that have reported a continuum of cell states as opposed to  
 442 sequential discrete cell types hierarchies<sup>48–51</sup>. In line with a recent murine study<sup>47</sup>, we have  
 443 shown that CD44 is expressed in naïve hPSCs-derived progenitors and here we  
 444 demonstrated that both human and mouse progenitors also express LMO4, a LIM-domain  
 445 protein<sup>52</sup>. Recent scRNAseq detected LMO4 in both human granulocyte progenitors in the  
 446 bone marrow<sup>48</sup> and adult mouse HSC<sup>53</sup>, but its associated proteins have not been  
 447 identified. We also reported high levels of *ID* genes within the progenitors, target genes of  
 448 BMP signaling known to be involved in HSC emergence<sup>54–56</sup>. IDs, like LMOs proteins, do  
 449 not present DNA binding domain and rather act through binding of other proteins in  
 450 complexes also involved in HSPC development<sup>57</sup> and erythropoiesis<sup>58</sup>. Overexpression of  
 451 ID2 in human HSC from cord blood has been reported to enhance their functional  
 452 stemness in vivo<sup>59</sup>, supporting the idea that this class of proteins might maintain the  
 453 progenitor status and thus might be useful in alternative programming strategies of hPSCs.  
 454 The use of scRNAseq on vast numbers of cells allows to detect even the rarest cell  
 455 population and we considered that it might enable the detection of rare HSC-like cells in  
 456 differentiating hPSCs cultures. We showed to hPSCs-derived cells showed a remarkable  
 457 specific expression pattern of marker genes identified in the human embryo, for example,  
 458 *SPINK2*, a novel marker of human fetal liver HSC and MPP, marked also our naïve  
 459 progenitors. By using machine learning we identify specific cell types sampled in vivo and  
 460 detected a small and transient population of HSC-like cells that, when compared to their in  
 461 vivo counterpart from fetal liver, showed only small transcriptional differences. Previous  
 462 reports described the hematopoietic progenitors obtained with the differentiation employed

463 in this work as intraembryonic-like<sup>18</sup>, using T-cells lineage as hallmark of definitive  
 464 hematopoiesis. However, yolk sac shows T-cell potential prior to HSC emergence<sup>60,61</sup>,  
 465 thus limiting the use of T-cell assay alone as discriminative of the corresponding  
 466 developmental wave. Our machine learning approach and the detection of HSC-like cells  
 467 strongly supports the intraembryonic identities of the hematopoietic cells differentiated in  
 468 vitro and provide an alternative and multifactorial approach to address questions regarding  
 469 the similarities to developmental counterparts. The unbiased and comprehensive  
 470 comparison used in this study allowed us to pinpoint differentially expressed genes  
 471 between in vitro-derived and in vivo HSCs can now be exploited to improve production of  
 472 HSC in vitro. Our analysis indicates in vitro HSC-like cells do not express CD43,  
 473 comparable to mouse Pro-HSC prior their maturation into functional definitive HSC<sup>62</sup>. This  
 474 could suggest that the widely acknowledged inability of hPSCs-derived progenitors to  
 475 reconstitute the hematopoietic system, could be due to their immature phenotype and the  
 476 lack of appropriate culture conditions for HSCs maturation and maintenance. In addition,  
 477 the identification of the HSC-like population using our machine learning approach, which is  
 478 based on high similarity in the gene expression profiles, could suggest that the molecular  
 479 basis of the functional deficiency of this in vitro derived population could reside at a post-  
 480 transcriptional level. Thus, future experiment will be required to assess whether a further  
 481 ad-hoc maturation step of sorted HSC-like cells would achieve reconstitution.  
 482 When we compared the hematopoietic progenitors developed in the human embryo  
 483 throughout gestations together with those derived in vitro, we found that while *SPINK2*  
 484 was expressed by all progenitors, *SPINK1* was exclusively detected in cells from the yolk  
 485 sac. *SPINK1* is able to bind to EGFR and induce epithelial to mesenchymal transition in  
 486 cancer cells<sup>63,64</sup>, a process similar to the endothelial to hematopoietic transition, where the  
 487 role of *SPINK1* remains largely unexplored. In summary, we propose here *SPINK1* as a  
 488 possible marker for primitive hematopoiesis which could be an extremely useful genes to  
 489 trace the cells that colonize the embryo from the yolk sac.  
 490 The differentiation protocol used in this study is well defined and serum-free and is one of  
 491 the most commonly used protocols used by many laboratories. We also showed that our  
 492 markers are able to identify functionally similar progenitors in different cell lines. Thus, our  
 493 browsable datasets and the findings of our study will be of interest to many in the field of  
 494 hematopoiesis and will allow to test how the frequency of this populations vary in response  
 495 to different cytokines conditions. In addition, the increasing availability of large scRNAseq

dataset of human tissue makes our pipeline applicable to the analyses of other systems where the hPSCs differentiation aims to produce adult-like cells for clinical application. In this way cell types differentiated in vitro can now be annotated in an unbiased manner that does not rely on few known markers and allows the identification of transcriptional discrepancies between cell types produced in vitro and their in vivo counterparts. In conclusion, our browsable dataset provides a comprehensive transcriptional characterization of in vitro derived hematopoietic progenitors. This work defines the makeup of the progenitor populations that give rise to immune cells, such as macrophages and granulocytes, as well as HSC-like cells, which holds great promise for their therapeutically application.

### **Acknowledgment**

The work was funded by Wellcome Trust (Grant No. 102610), MRC Innovate UK (Grant No. 102853), BBSRC (Grant No. S002219/1). AF received a Carnegie Incentive Grant (Grant No. RIG008218). AB received BBSRC Future Leaders Fellowship (Grant reference BB/N011597/1). Sequencing was carried out by Edinburgh Genomics, The University of Edinburgh. Edinburgh Genomics is partly supported through core grants from NERC (R8/H10/56), MRC (MR/K001744/1) and BBSRC (BB/J004243/1). We thank Professor Ben Macarthur for suggesting the application of machine learning to our study; Andrew Elefanty for sharing the RUNX1C-GFP cell line; Fiona Rossi, Claire Cryer, Bindi Heer and Andrea Corsinotti from the Flow Facility as well as Bertand Verney and Matthieu Vermeren from the imaging facility. This work has made use of the resources provided by the Edinburgh Compute and Data Facility (ECDF) (<http://www.ecdf.ed.ac.uk/>).

### **Author Contribution**

AF, designed and performed research, analyzed the data and wrote the manuscript. AF, PS, AB and NR performed bioinformatics analysis. PR, ST, MLY, AHT, JE, BH, RA performed research. LMF designed the experiment, analyzed data and wrote the manuscript. NH, AM, KO, and NR provided intellectual input and final approval of the manuscript.

### **Declaration of interest**

Authors declare no competing interests.

529

530 **References**

531

- 532 1. Vo LT, Daley GQ. De novo generation of HSCs from somatic and pluripotent stem  
533 cell sources. *Blood*. 2015;125(17):2641–8.
- 534 2. Palis J. Hematopoietic stem cell-independent hematopoiesis: emergence of  
535 erythroid, megakaryocyte, and myeloid potential in the mammalian embryo. *FEBS*  
536 *Lett*. 2016;590(22):3965–3974.
- 537 3. Tober J, Koniski A, McGrath KE, et al. The megakaryocyte lineage originates from  
538 hemangioblast precursors and is an integral component both of primitive and of  
539 definitive hematopoiesis. *Blood*. 2007;109(4):1433–41.
- 540 4. Palis J, Robertson S, Kennedy M, Wall C, Keller G. Development of erythroid and  
541 myeloid progenitors in the yolk sac and embryo proper of the mouse. *Development*.  
542 1999;126(22):5073–84.
- 543 5. McGrath KE, Frame JM, Fegan KH, et al. Distinct Sources of Hematopoietic  
544 Progenitors Emerge before HSCs and Provide Functional Blood Cells in the  
545 Mammalian Embryo. *Cell Rep*. 2015;
- 546 6. Mass E, Ballesteros I, Farlik M, et al. Specification of tissue-resident macrophages  
547 during organogenesis. *Science (80-. )*. 2016;353:6304.
- 548 7. Schulz C, Gomez Perdiguero E, Chorro L, et al. A lineage of myeloid cells  
549 independent of Myb and hematopoietic stem cells. *Science (80-. )*.  
550 2012;336(6077):86–90.
- 551 8. Stremmel C, Schuchert R, Wagner F, et al. Yolk sac macrophage progenitors traffic  
552 to the embryo during defined stages of development. *Nat. Commun*. 2018;9(1):75.
- 553 9. Medvinsky A, Dzierzak E. Definitive hematopoiesis is autonomously initiated by the  
554 AGM region. *Cell*. 1996;86(6):897–906.
- 555 10. Ivanovs A, Rybtsov S, Welch L, et al. Highly potent human hematopoietic stem cells  
556 first emerge in the intraembryonic aorta-gonad-mesonephros region. *J. Exp. Med*.  
557 2011;208(12):2417–2427.
- 558 11. Easterbrook J, Fidanza A, Forrester LM. Concise review: Programming human  
559 pluripotent stem cells into blood. *Br. J. Haematol*. 2016;173(5):.
- 560 12. Ivanovs A, Rybtsov S, Anderson RA, Turner ML, Medvinsky A. Identification of the  
561 niche and phenotype of the first human hematopoietic stem cells. *Stem Cell Reports*.  
562 2014;2(4):449–456.
- 563 13. Tavian M, Hallais MF, Péault B. Emergence of intraembryonic hematopoietic  
564 precursors in the pre-liver human embryo. *Development*. 1999;126(4):793–803.
- 565 14. Guibentif C, Rönn RE, Böiers C, et al. Single-Cell Analysis Identifies Distinct Stages  
566 of Human Endothelial-to-Hematopoietic Transition. *Cell Rep*. 2017;19(1):10–19.
- 567 15. Angelos MG, Abrahante JE, Blum RH, Kaufman DS. Single Cell Resolution of  
568 Human Hematoendothelial Cells Defines Transcriptional Signatures of Hemogenic  
569 Endothelium. *Stem Cells*. 2018;36(2):206–217.
- 570 16. Han X, Chen H, Huang D, et al. Mapping human pluripotent stem cell differentiation  
571 pathways using high throughput single-cell RNA-sequencing. *Genome Biol*.  
572 2018;19(1):1–19.
- 573 17. Stoeckius M, Hafemeister C, Stephenson W, et al. Simultaneous epitope and  
574 transcriptome measurement in single cells. *Nat. Methods*. 2017;14(9):865–868.
- 575 18. Sturgeon CM, Ditadi A, Awong G, Kennedy M, Keller G. Wnt signaling controls the  
576 specification of definitive and primitive hematopoiesis from human pluripotent stem

- cells. *Nat. Biotechnol.* 2014;32(6):554–561.
19. Vodyanik MA, Thomson JA, Slukvin II, Dulac C, Péault B. Leukosialin (CD43) defines hematopoietic progenitors in human embryonic stem cell differentiation cultures. *Blood.* 2006;108(6):2095–105.
  20. Garcia-Alegria E, Menegatti S, Fadlullah MZH, et al. Early Human Hemogenic Endothelium Generates Primitive and Definitive Hematopoiesis In Vitro. *Stem Cell Reports.* 2018;11(5):1061–1074.
  21. Ng ES, Azzola L, Bruveris FF, et al. Differentiation of human embryonic stem cells to HOXA+ hemogenic vasculature that resembles the aorta-gonad-mesonephros. *Nat. Biotechnol.* 2016;34(11):1168–1179.
  22. Sroczynska P, Lancrin C, Kouskoff V, Lacaud G. The differential activities of Runx1 promoters define milestones during embryonic hematopoiesis. *Blood.* 2009;114(26):5279–89.
  23. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* 2018;36(5):411–420.
  24. Haghverdi L, Büttner M, Wolf FA, Buettner F, Theis FJ. Diffusion pseudotime robustly reconstructs lineage branching. *Nat. Methods.* 2016;13(10):845–848.
  25. Qiu X, Mao Q, Tang Y, et al. Reversed graph embedding resolves complex single-cell trajectories. *Nat. Methods.* 2017;14(10):979–982.
  26. Wolf FA, Hamey FK, Plass M, et al. PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome Biol.* 2019;20(1):1–9.
  27. Moreau T, Evans AL, Vasquez L, et al. Large-scale production of megakaryocytes from human pluripotent stem cells by chemically defined forward programming. *Nat. Commun.* 2016;7(1):11208.
  28. Stumpf PS, Smith RCG, Lenz M, et al. Stem Cell Differentiation as a Non-Markov Stochastic Process. *Cell Syst.* 2017;5(3):268–282.e7.
  29. Chan TE, Stumpf MPH, Babbie AC. Gene Regulatory Network Inference from Single-Cell Data Using Multivariate Information Measures. *Cell Syst.* 2017;5(3):251–267.e3.
  30. Lopez-Yrigoyen M, Fidanza A, Cassetta L, et al. A human iPSC line capable of differentiating into functional macrophages expressing ZsGreen: A tool for the study and in vivo tracking of therapeutic cells. *Philos. Trans. R. Soc. B Biol. Sci.* 2018;373(1750):.
  31. Jaffredo T, Gautier R, Eichmann A, Dieterlen-Lièvre F. Intraaortic hemopoietic cells are derived from endothelial cells during ontogeny. *Development.* 1998;125(22):4575–83.
  32. Zovein AC, Hofmann JJ, Lynch M, et al. Fate Tracing Reveals the Endothelial Origin of Hematopoietic Stem Cells. *Cell Stem Cell.* 2008;3(6):625–636.
  33. Jayapal SR, Lee KL, Ji P, et al. Down-regulation of Myc is essential for terminal erythroid maturation. *J. Biol. Chem.* 2010;285(51):40252–65.
  34. Lammers R, Giesert C, Grünebach F, et al. Monoclonal antibody 9C4 recognizes epithelial cellular adhesion molecule, a cell surface antigen expressed in early steps of erythropoiesis. *Exp. Hematol.* 2002;30(6):537–545.
  35. Clay D, Rubinstein E, Mishal Z, et al. CD9 and megakaryocyte differentiation. *Blood.* 2001;97(7):1982–1989.
  36. Kalev-Zylinska ML, Horsfield JA, Flores MVC, et al. Runx3 is required for hematopoietic development in zebrafish. *Dev. Dyn.* 2003;228(3):323–336.
  37. Fainaru O, Woolf E, Lotem J, et al. Runx3 regulates mouse TGF- $\beta$ -mediated

- dendritic cell function and its absence results in airway inflammation. *EMBO J.* 2004;23:969–979.
38. Popescu DM, Botting RA, Stephenson E, et al. Decoding human fetal liver haematopoiesis. *Nature.* 2019;574(7778):365–371.
39. Stumpf PS, Du D, Imanishi H, et al. Mapping biology from mouse to man using transfer learning. *bioRxiv.* 2019;
40. Zeng Y, He J, Bai Z, et al. Tracing the first hematopoietic stem cell generation in human embryo by single-cell RNA sequencing. *Cell Res.* 2019;29(11):881–894.
41. Dunwoodie SL. The Role of Hypoxia in Development of the Mammalian Embryo. *Dev. Cell.* 2009;17(6):755–773.
42. Imanirad P, Solaimani Kartalaei P, Crisan M, et al. HIF1 $\alpha$  is a regulator of hematopoietic progenitor and stem cell development in hypoxic sites of the mouse embryo. *Stem Cell Res.* 2014;12(1):24–35.
43. Bigas A, Espinosa L. Hematopoietic stem cells: To be or Notch to be. *Blood.* 2012;119(14):3226–3235.
44. Wölwer CB, Pase LB, Russell SM, Humbert PO. Calcium signaling is required for erythroid enucleation. *PLoS One.* 2016;11(1):.
45. Lopez-Yrigoyen M, Yang C-T, Fidanza A, et al. Genetic programming of macrophages generates an in vitro model for the human erythroid island niche. *Nat. Commun.* 2019;10(1):881.
46. Yang C-T, Ma R, Axton RA, et al. Activation of KLF1 Enhances the Differentiation and Maturation of Red Blood Cells from Human Pluripotent Stem Cells. *Stem Cells.* 2017;35(4):886–897.
47. Oatley M, Bölükbaşı ÖV, Svensson V, et al. Single-cell transcriptomics identifies CD44 as a marker and regulator of endothelial to haematopoietic transition. *Nat. Commun.* 2020;11(1):.
48. Paul F, Arkin Y, Giladi A, et al. Transcriptional Heterogeneity and Lineage Commitment in Myeloid Progenitors. *Cell.* 2015;163(7):1663–1677.
49. Velten L, Haas SF, Raffel S, et al. Human haematopoietic stem cell lineage commitment is a continuous process. *Nat. Cell Biol.* 2017;19(4), 271–281.
50. Rodriguez-Fraticelli AE, Wolock SL, Weinreb CS, et al. Clonal analysis of lineage fate in native haematopoiesis. *Nature.* 2018;553(7687):212–216.
51. Knapp DJHF, Hammond CA, Wang F, et al. A topological view of human CD34+ cell state trajectories from integrated single-cell output and proteomic data. *Blood.* 2019;133(9):927–939.
52. Grutz G, Forster A, Rabbitts TH. Identification of the LMO4 gene encoding an interaction partner of the LIM-binding protein LDB1/NLI1: a candidate for displacement by LMO proteins in T cell acute leukaemia. *Oncogene.* 1998;17(21):2799–2803.
53. Lai S, Huang W, Xu Y, et al. Cell Discovery Comparative transcriptomic analysis of hematopoietic system between human and mouse by Microwell-seq. *Cell Discov.* 2018;4:34.
54. Souilhol C, Gonneau C, Lendinez JG, et al. Inductive interactions mediated by interplay of asymmetric signalling underlie development of adult haematopoietic stem cells. *Nat. Commun.* 2016;(2016): 1–13.
55. Crisan M, Kartalaei PS, Vink CS, et al. BMP signalling differentially regulates distinct haematopoietic stem cell types. *Nat. Commun.* 2015;6(1):8040.
56. McGarvey AC, Rybtsov S, Souilhol C, et al. A molecular roadmap of the AGM region reveals BMPER as a novel regulator of HSC maturation. *J. Exp. Med.*

2017;214(12):3731–3751.

57. Wilson NK, Foster SD, Wang X, et al. Combinatorial Transcriptional Control In Blood Stem/Progenitor Cells: Genome-wide Analysis of Ten Major Transcriptional Regulators. *Cell Stem Cell*. 2010;7(4):532–544.
58. Wadman IA, Osada H, Grütz GG, et al. The LIM-only protein Lmo2 is a bridging molecule assembling an erythroid, DNA-binding complex which includes the TAL1, E47, GATA-1 and Ldb1/NLI proteins. *EMBO J*. 1997;16(11):3145–57.
59. van Galen P, Kreso A, Wienholds E, et al. Reduced Lymphoid Lineage Priming Promotes Human Hematopoietic Stem Cell Expansion. *Cell Stem Cell*. 2014;14(1):94–106.
60. Yoshimoto M, Porayette P, Glosson NL, et al. Autonomous murine T-cell progenitor production in the extra-embryonic yolk sac before HSC emergence. *Blood*. 2012;119(24):5706–14.
61. Gentek R, Ghigo C, Hoeffel G, et al. Epidermal  $\gamma\delta$  T cells originate from yolk sac hematopoiesis and clonally self-renew in the adult. *J. Exp. Med*. 2018;215(12):2994–3005.
62. Rybtsov S, Batsivari A, Bilotkach K, et al. Tracing the origin of the HSC hierarchy reveals an SCF-dependent, IL-3-independent CD43- embryonic precursor. *Stem Cell Reports*. 2014;3(3):489–501.
63. Wang C, Wang L, Su B, et al. Serine protease inhibitor Kazal type 1 promotes epithelial-mesenchymal transition through EGFR signaling pathway in prostate cancer. *Prostate*. 2014;74(7):689–701.
64. Chen F, Long Q, Fu D, et al. Targeting SPINK1 in the damaged tumour microenvironment alleviates therapeutic resistance. *Nat. Commun*. 2018;9(1):1–19.

## Figure Legends

### Figure 1 - Single cell transcriptome analysis reveals clusters of naïve and lineage-committed haematopoietic progenitors.

(A) Schematic of the single cell RNA sequencing experiment where iPSCs (SFCi55) were differentiated in vitro (IVD) for 13 days (Supplementary Figure1A), CD235a<sup>+</sup>CD43<sup>+</sup> suspension cells were isolated by flow cytometry and subjected to 10x genomics sequencing platform. (B) tSNE visualization of 11,420 cells divided into 8 clusters including clusters defined by gene expression as naïve (naïve 1 and 2), and others that expressed genes associated with erythroid (Ery1 and 2), megakaryocyte (Mega 1,2 and 3) and granulocyte (granulo) lineages. (C) Heatmap showing expression of the top 10 marker genes for each cluster (colors for each cluster as in Figure 1B). (D) Gene expression levels of marker genes associated with different progenitor cell types that were identified by clustering, visualized on tSNE.

### Figure 2 - Trajectory analyses support naïve progenitor identity and their progression to lineage committed progenitors.

(A) Diffusion plot displays the naïve progenitors in the core region of the plot from where the three direction of commitment originates, the arrows indicates the commitment directions. (B) Representation of each cluster on the diffusion plot. (C) PAGA analysis show that naïve cluster are connected to the lineage committed cells. Each node contains



a pie chart showing the proportion of cells for each cluster. Colors indicate cluster identities. **(D)** Monocle trajectory analyses demonstrates a similar pattern to that obtained from the diffusion plot shown in A with naïve progenitors at the top of the hierarchy, with progression toward committed. **(E)** Monocle trajectory visualizing each cluster individually. **(F)** Expression levels of the marker genes coding for transcription factors associated with each cluster, bars color indicates the cluster.

### Figure 3 - CD44 identifies clonogenic hematopoietic progenitors.

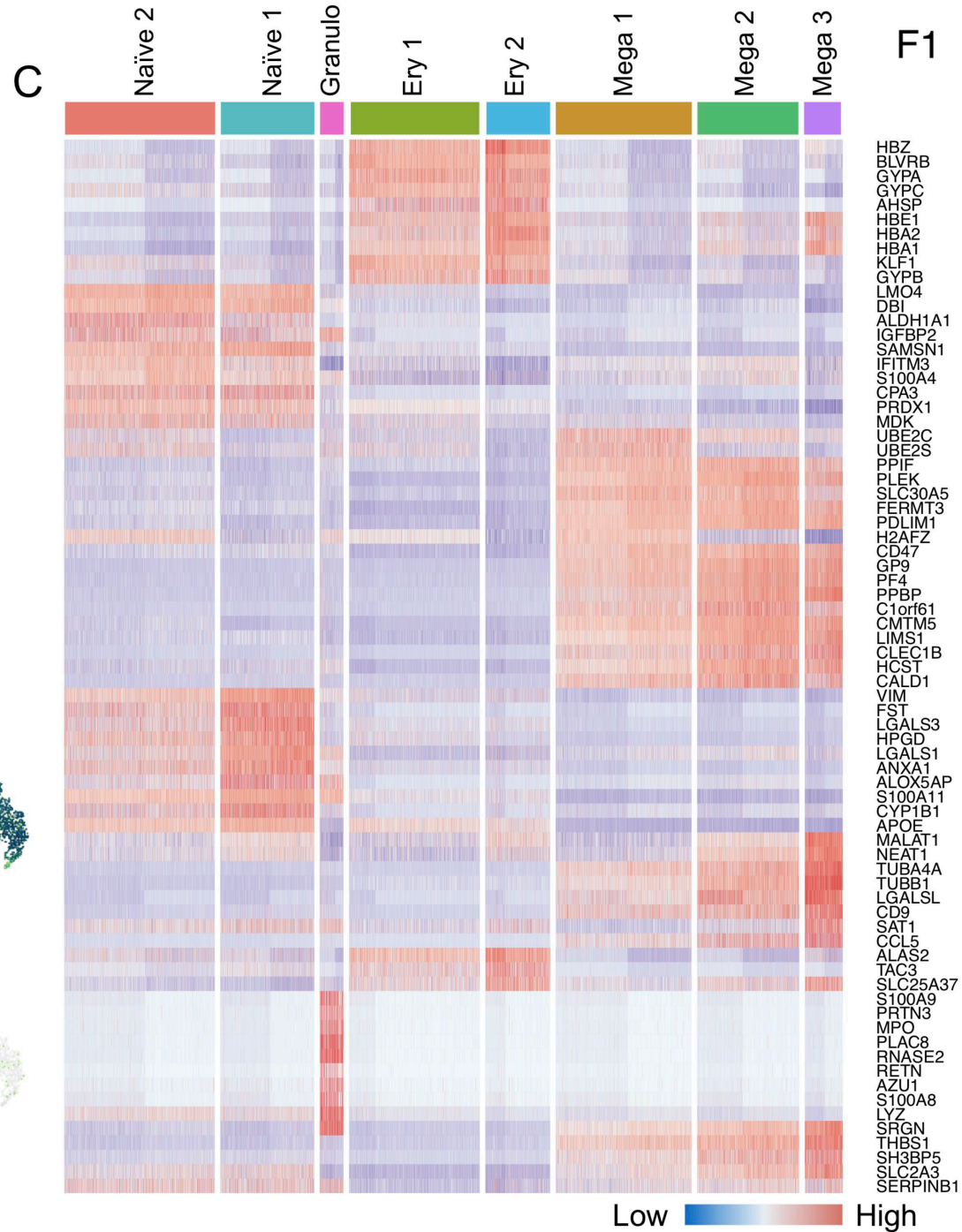
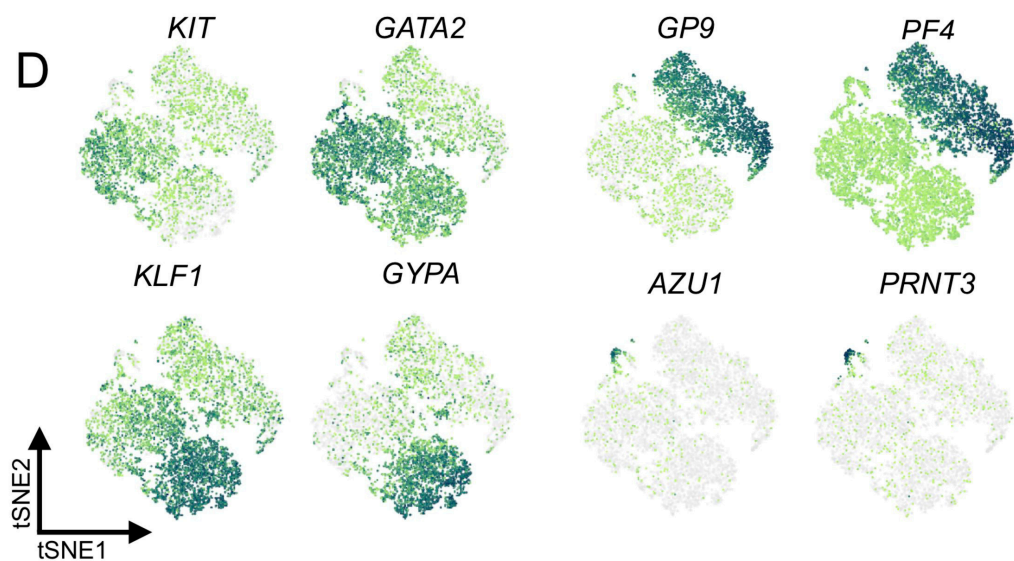
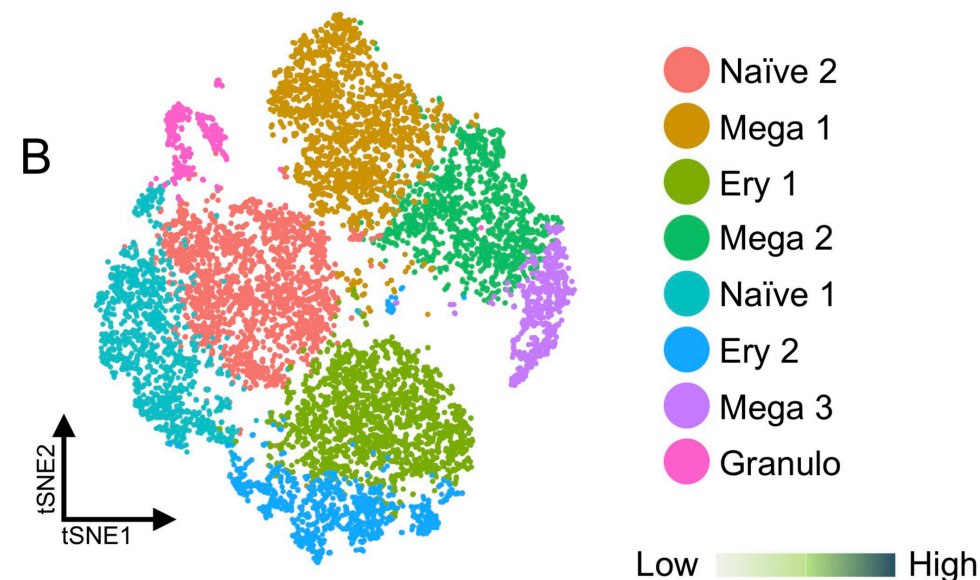
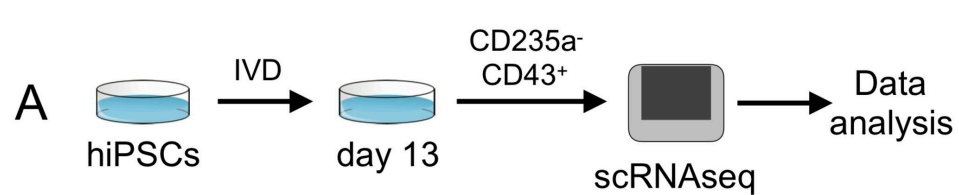
**(A)** Expression levels of genes encoding cell surface markers, *CD33*, *CD44* and *CD18* that were associated with the naïve progenitor clusters. **(B)** Scatter plot of flow cytometry profile of naïve 1A, 1B and 2 cells at day 13 of differentiation (hiPSCs-SFCi55). Cells are gated on *CD235a*<sup>-</sup>*CD43*<sup>+</sup>*CD33*<sup>+</sup>. **(C)** Schematic of the chimeric culture system using the ZsGreen reporter to trace cells during the differentiation process. ZsGreen and parental line (SFCi55) were differentiated in a synchronous manner, at day 10 naïve 1 and naïve 2 cells are sorted and co-cultured with the parental line differentiation. Co-culture is then analyzed at day 13. **(D)** Representative flow cytometry profile of the day 13 naïve descendants' cells after sorting at day 10 and chimeric co-culturing of naïve 1 (teal) and naïve 2 (pink) cells. Contribution of naïve 1, in teal, and naïve 2, in pink, to the naïve 1A, 1B and 2 compartment (n=6, multinomial logistic regression, \*p<0.05, \*\*p<0.01, \*\*\*p<0.005). **(E)** CFU-C analyses of FAC-sorted naïve 1 and naïve 2 cells from day 10 (n=3, paired t-Test p=0.0753) (hiPSCs-SFCi55). **(F)** CFU-C analyses of FAC-sorted naïve 1A, 1B and 2 cells from day 13 (n=9, Holm-Sidak's test, p<=0.001) (hiPSCs-SFCi55). **(G)** RUNX1-GFP expression in naïve 1 and naïve 2 at both day 10 and day 13 (n=12, paired t test; \* p<0.05, \*\* p<0.01).

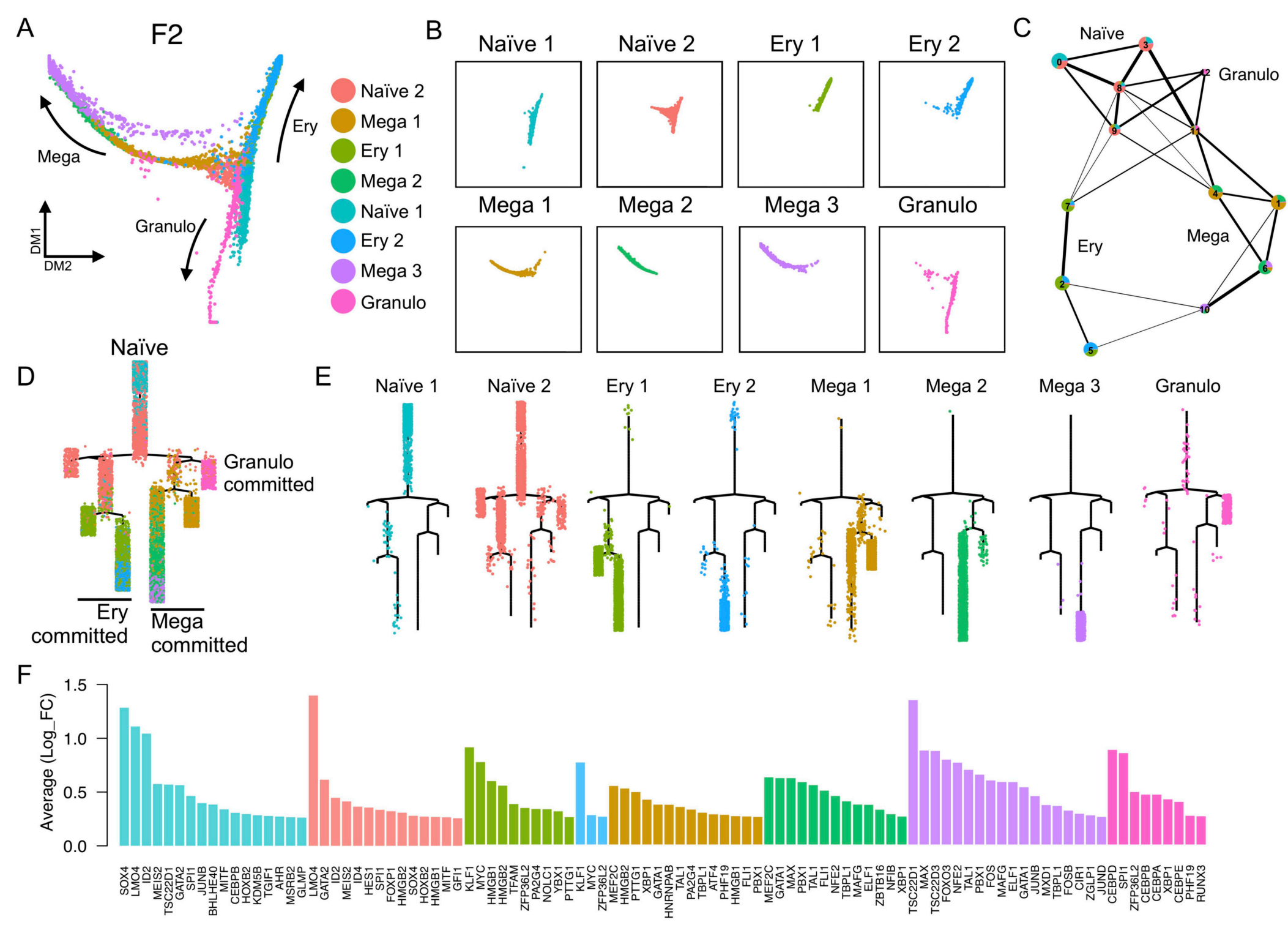
### Figure 4 - CITE-seq analyses confirm markers for naïve and lineage-committed progenitor cells.

**(A)** tSNE visualization of the CITE-seq analysis of day 10 and day 13 *CD235a*<sup>-</sup> suspension cells; reduction and clustering were performed using only transcriptomic data (hiPSC-SFCi55). **(B)** tSNE visualization of the libraries obtained from *CD235a*<sup>-</sup> suspension cells collected at day 10 (pink) and day 13 (teal) showing lineage commitment direction. **(C)** Single cell protein expression level of the membrane markers associated with the different cell types. Data are visualized on tSNE (ADT = antibody derived tags). **(D)** tSNE plot and annotation of clustering obtained from analysis derived from ADT data alone. **(E)** Visualization of the libraries obtained from *CD235a*<sup>-</sup> suspension cells projected on the tSNE obtained from ADT data in D, cell progression shows lineage commitment trajectory. Cells are colored according to the day of collection (day 10 in pink and day 13 in teal). **(F)** Confusion matrix of clustering obtained from complete transcriptomic data (RNA) and that obtained from ADT data alone. Color of each box indicates the % of cells classified in each RNA versus ADT cluster. **(G)** Gene expression levels of human foetal gene marker genes in our naïve and lineage committed progenitors.

**Figure 5 - Artificial neural network identifies HSC-like cells in iPSC derived hematopoietic cells.**

**(A)** Schematic of the artificial neural network (ANN) architecture for label-transfer. An input layer (3479 units), two fully connected hidden layers (64 and 32 units) and a 28-unit softmax layer corresponding to cell types in the source domain (human foetal liver scRNAseq data) used for training. Classification of cell types in the target domain of human iPSC-derived single cell transcriptomes (test data). **(B)** Proportion of cells labelled HSC/MPPs at day 10 or day 13 of hiPSC differentiation in vitro. **(C)** Proportion of in vitro derived CD235a<sup>-</sup> progenitors and CD235a<sup>-</sup>CD43<sup>+</sup> cells labelled HSC/MPP by the ANN (ND = not detected). **(D)** UMAP visualization of the integrated dataset containing in vivo derived (blue) and in vitro annotated (pink) HSC/MPPs. **(E)** Heatmap of differentially expressed genes coding for transcription factors obtained from the comparison of in vivo and in vitro derived HSC/MPPs. **(F)** Transcription factor binding motifs enriched upstream of the 54 genes identified as differentially expressed between HSC/MPP generated in vitro and in vivo.

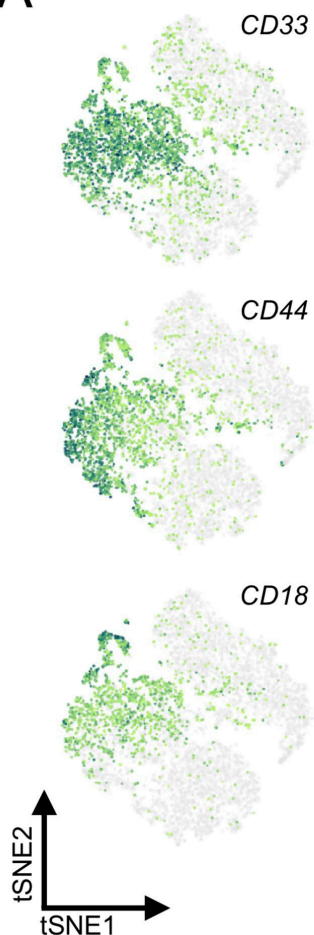




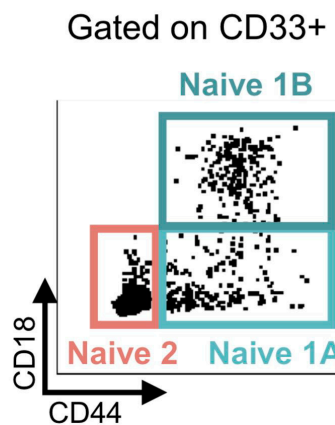


F3

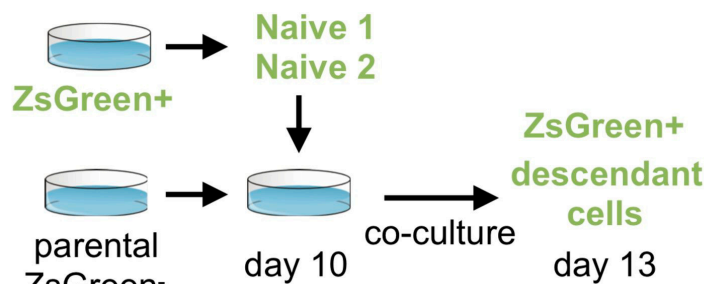
A



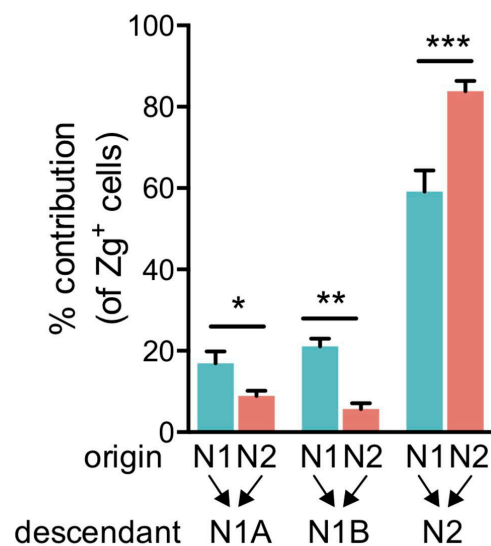
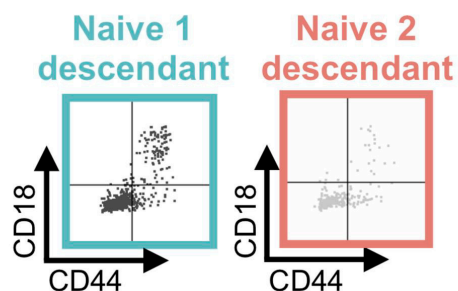
B



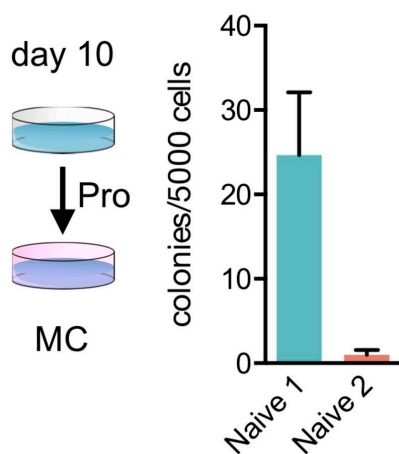
C



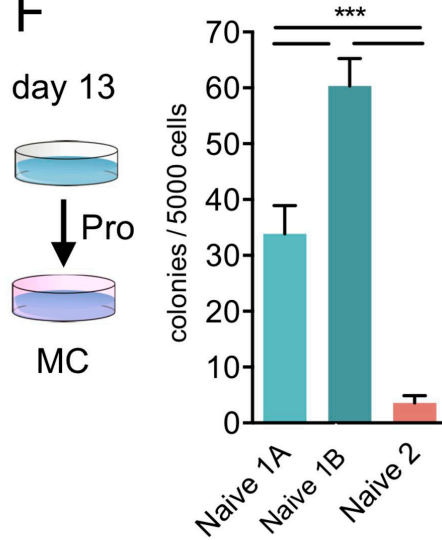
D



E



F



G

